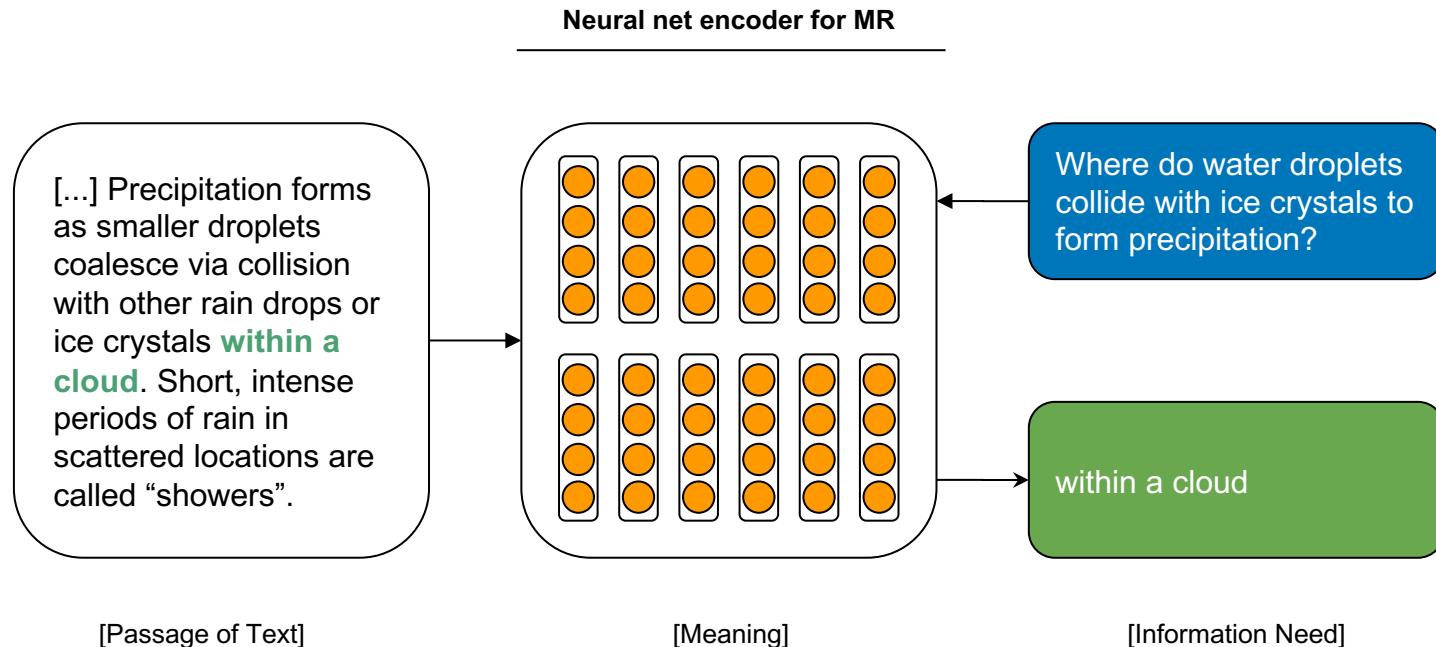


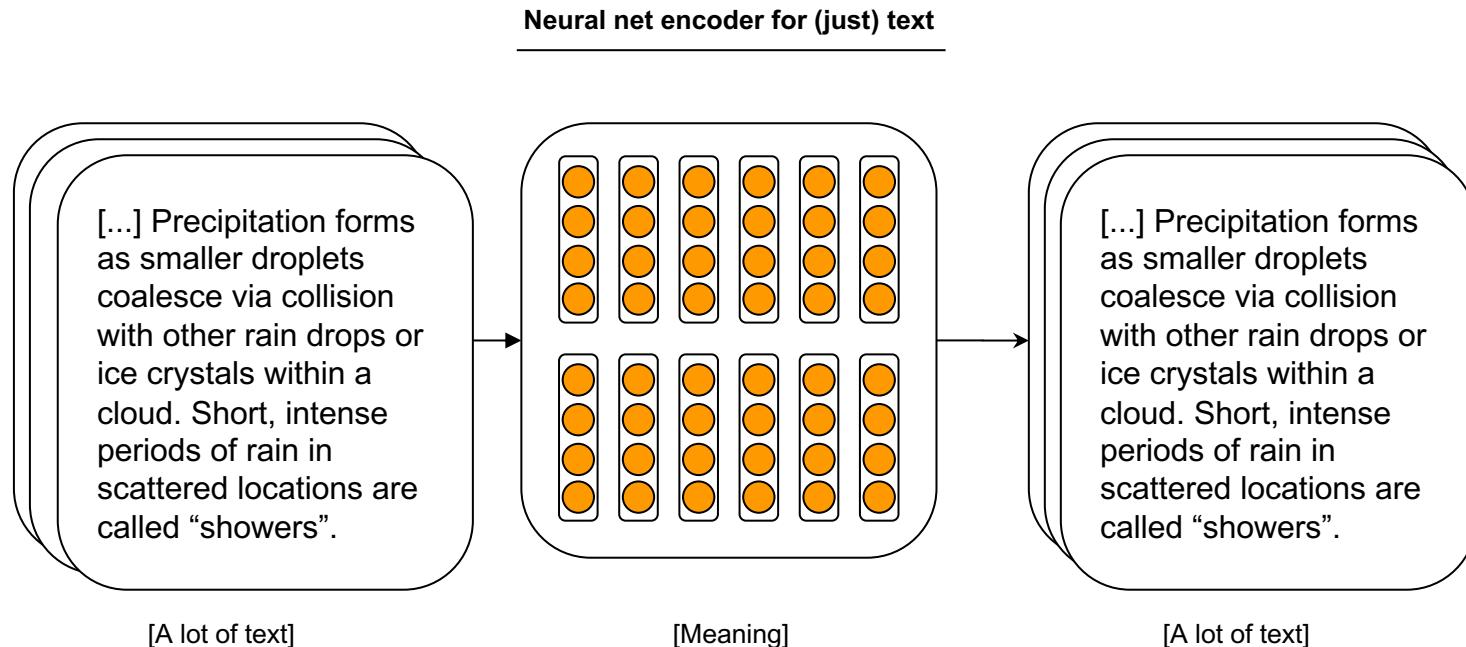
# Machine Reading / Current Trend

---

# Supervised training

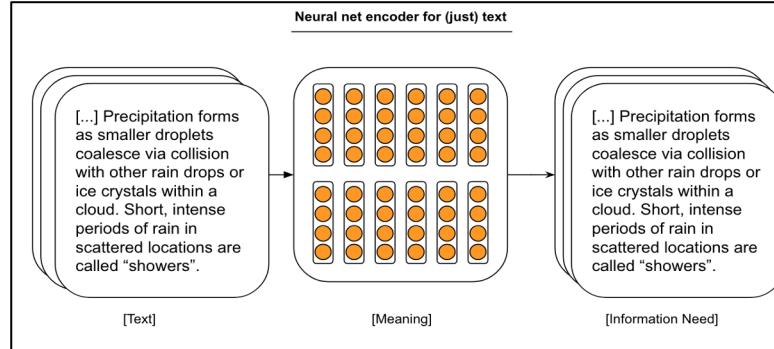


# Unsupervised pretrained representations

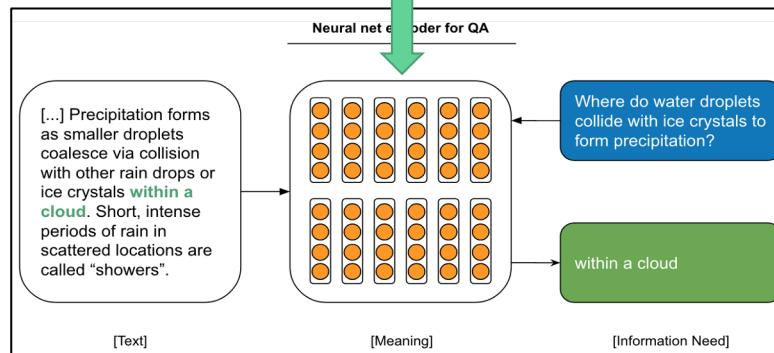


# Lifting over pretrained representations

## Pretrained Language Model



Transfer



# How is this different from pretrained word embeddings?

## Pretrained **Word** Embeddings (word2vec)

- Predicting co-occurring of words
- Independent of other context

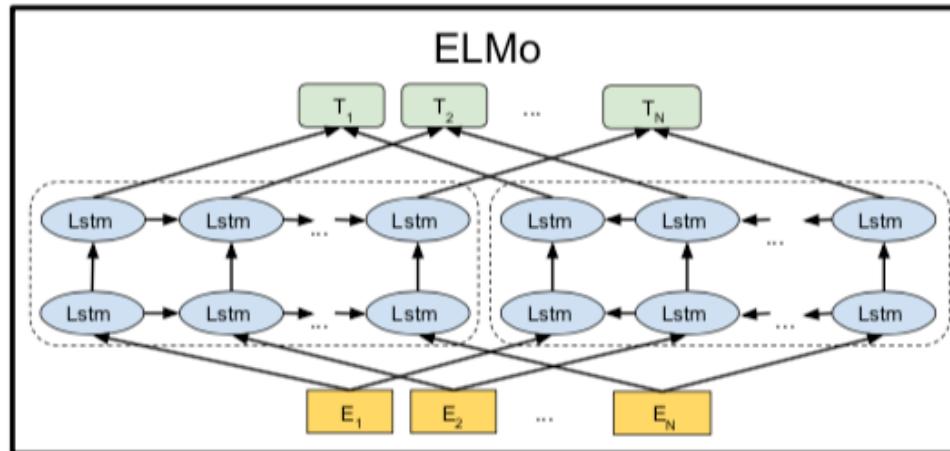
## Pretrained **Contextualized** Embeddings (e.g. ELMo, BERT)

- Predicting whole text (using LSTM, or Self-Attention)
- Full dependence on other context

# ELMo: Embeddings from Language Models

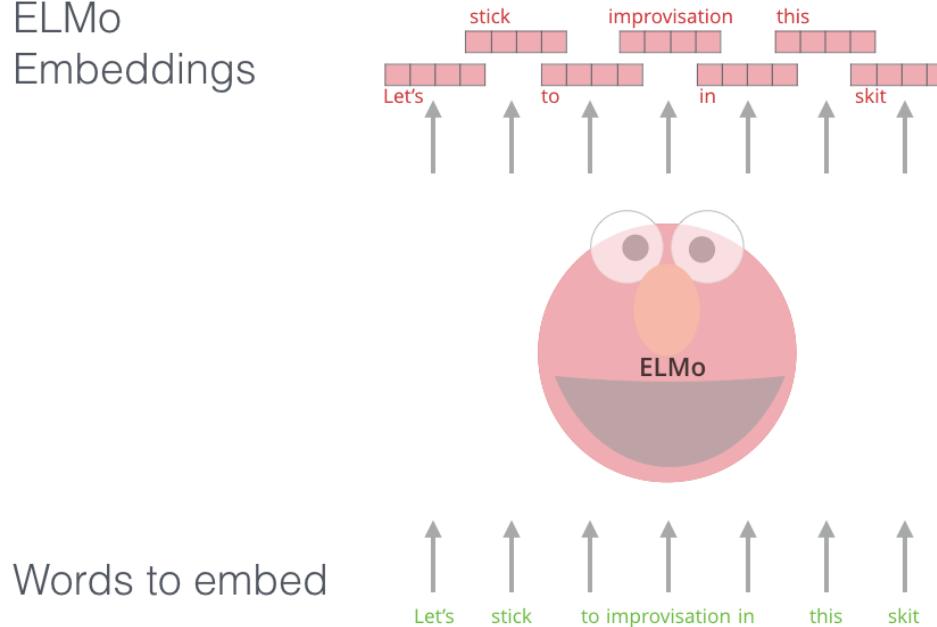
Peters et al., NAACL'18

- Train a BiLSTM for Bidirectional language modeling on a large dataset
- Run the sentence to encode through both forward and backward LSTMs
- Combine forward and backward representations into final contextual embeddings



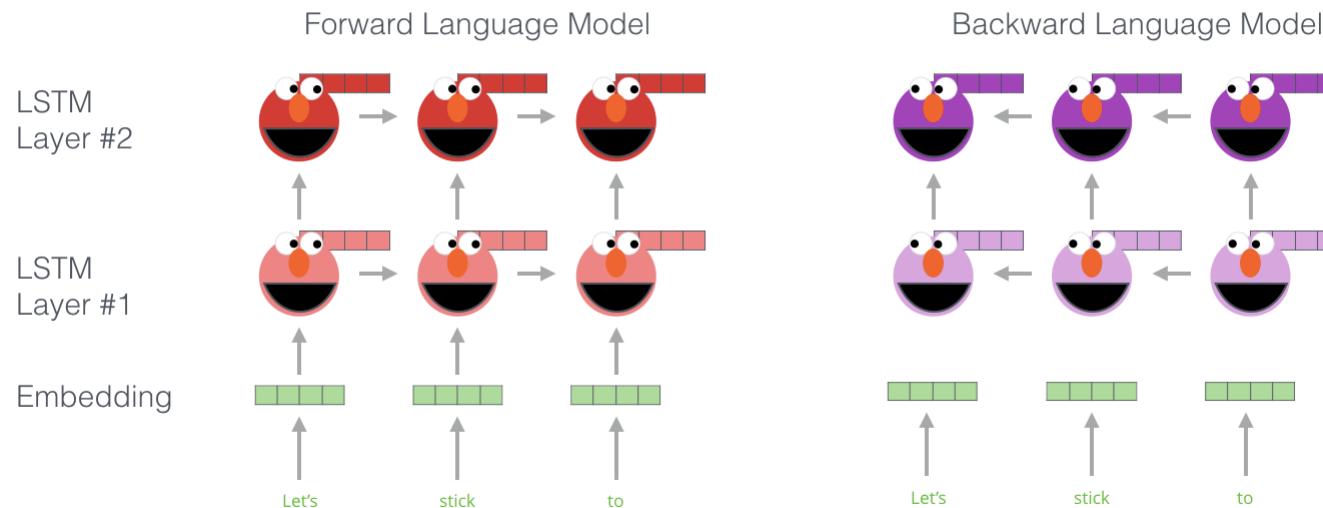
# ELMo: Embeddings from Language Models

ELMo  
Embeddings



# ELMo: Embeddings from Language Models

Embedding of “stick” in “Let’s stick to” - Step #1



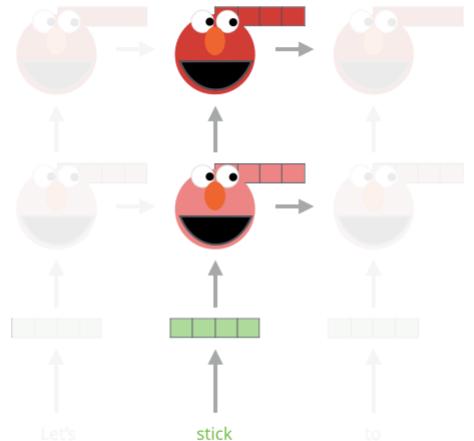
# ELMo: Embeddings from Language Models

Embedding of “stick” in “Let’s stick to” - Step #2

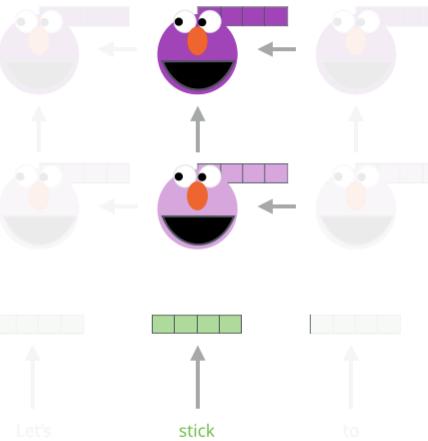
1- Concatenate hidden layers



Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task

$$\text{Red vector} \times s_2$$

$$\text{Purple vector} \times s_1$$

$$\text{Green vector} \times s_0$$

3- Sum the (now weighted) vectors



ELMo embedding of “stick” for this task in this context

# ELMo performance

Task	Previous SOTA		Our Baseline	ELMo + Baseline	Increase (Absolute/Relative)
Machine Reading - SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
Textual Entailment - SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
Semantic Labeling - SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coreference Resolution - Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
Entity Extraction - NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
Sentiment Analysis - SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

# What is ELMo learning ?

- Meaning of words in context
  - POS, word sense, etc.

	Source	Nearest Neighbors
GloVe	play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM	Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent <u>play</u> .
	Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}	{...} they were actors who had been handed fat roles in a successful <u>play</u> , and had talent enough to fill the roles competently , with nice understatement .

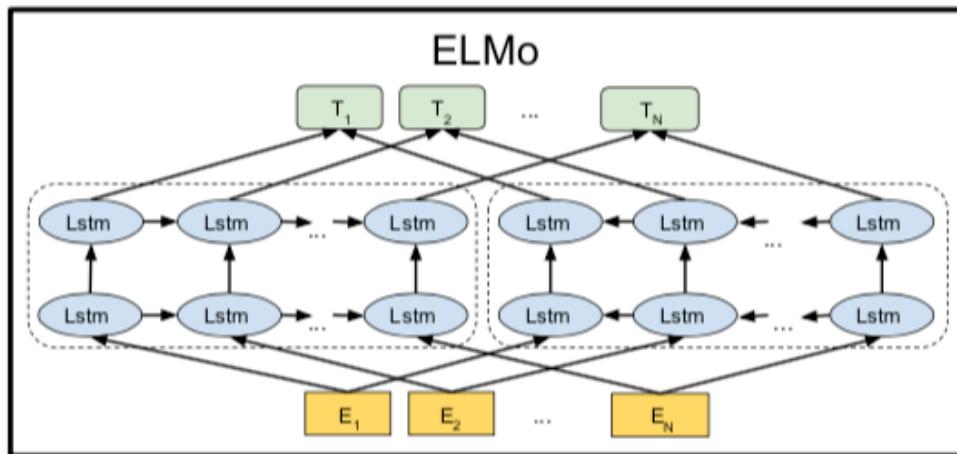
# Problems with ELMo

- Need to use different architectures for different tasks
- Retraining models is slow, transfer learning is fast
- Need to deal with long term dependencies in LSTMs!

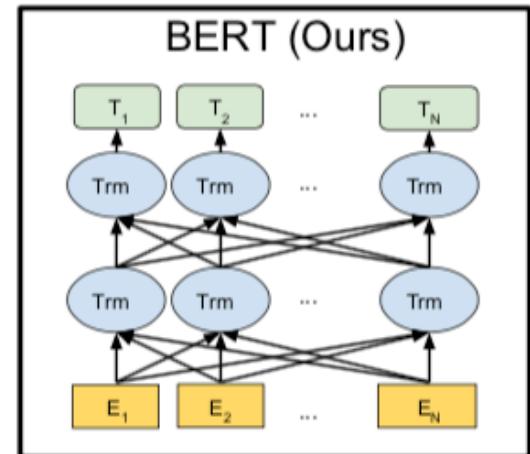
# BERT - Bidirectional Encoder Representations from Transformers

Devlin et al., NAACL'19

Solutions: use Transformer + encoder layers instead of decoder layers



(OpenAI GPT)



Innovation with multiple pretraining tasks

# BERT – Pretraining 1: masked language modeling

- Given a sentence with some words masked at random, can we predict them?
- Randomly select 15% of tokens to be replaced with “<MASK>”

# BERT – Pretraining 1: masked language modeling

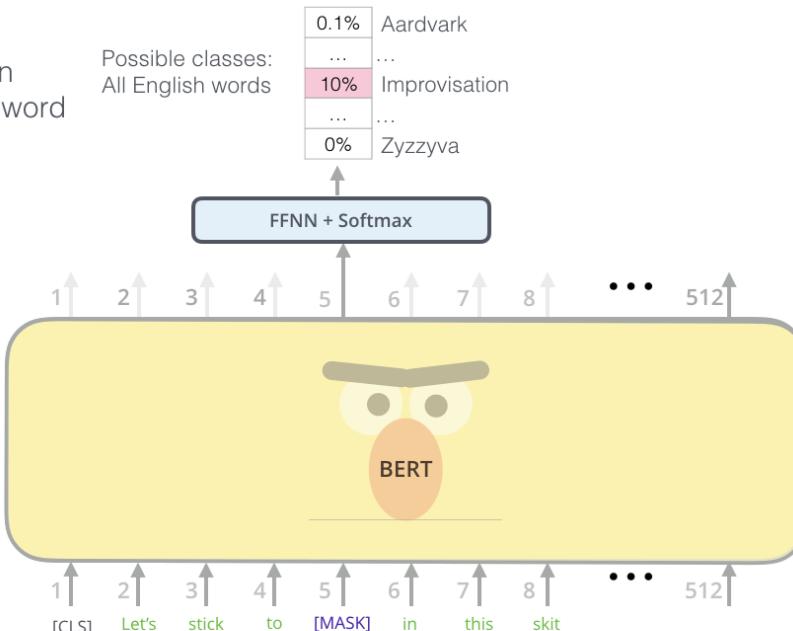
Use the output of the masked word's position to predict the masked word

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax

Randomly mask  
15% of tokens



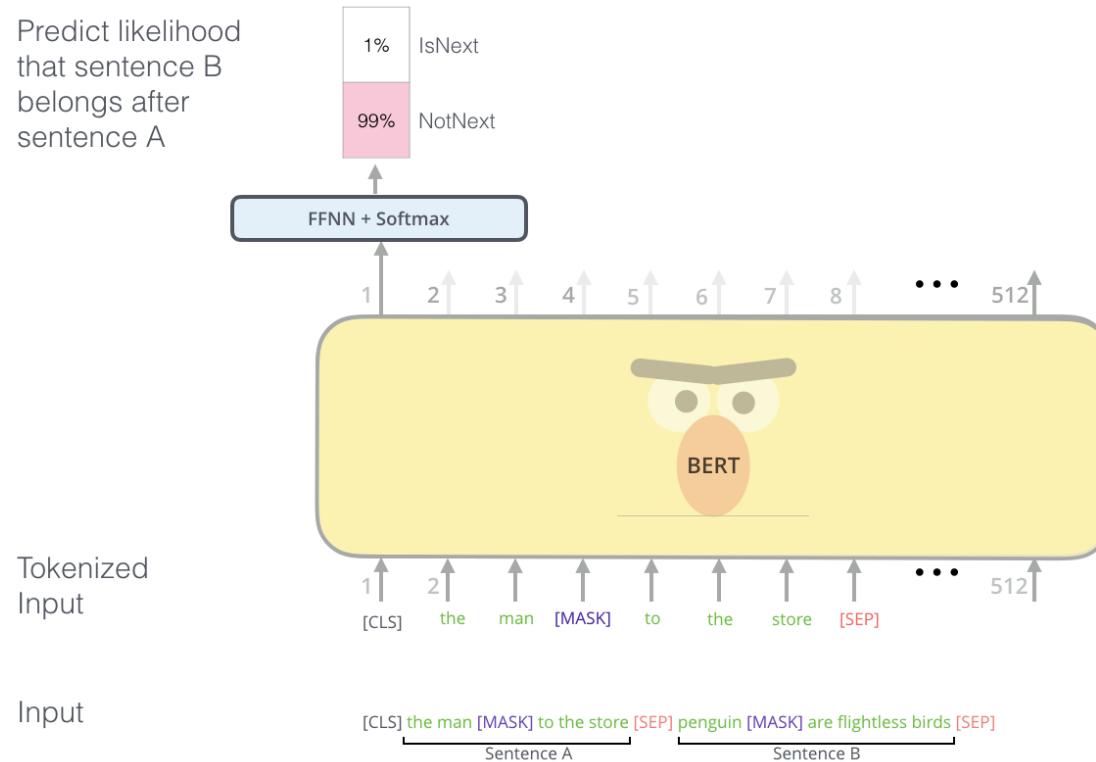
Input

[CLS] Let's stick to improvisation in this skit

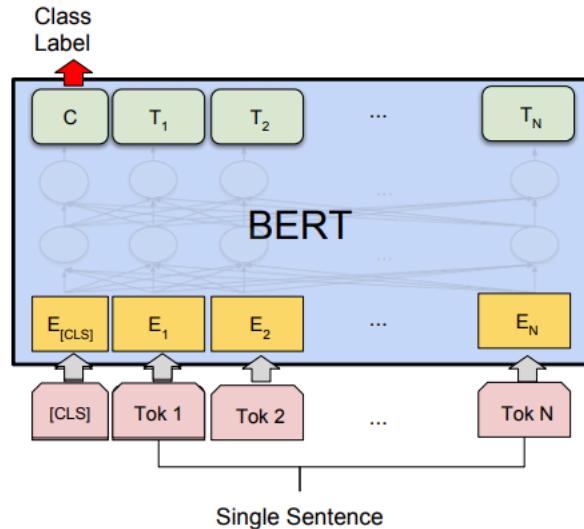
## BERT – Pretraining 2: next sentence prediction

- Given two sentences, does the first follow the second?
- Teaches BERT about relationship between two sentences
- 50% of the time the actual next sentence, 50% random

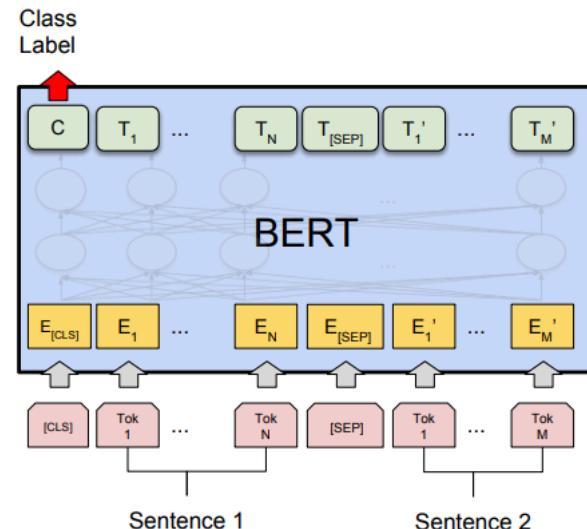
# BERT – Pretraining 2: next sentence prediction



# BERT – Fine-tuning for Classification

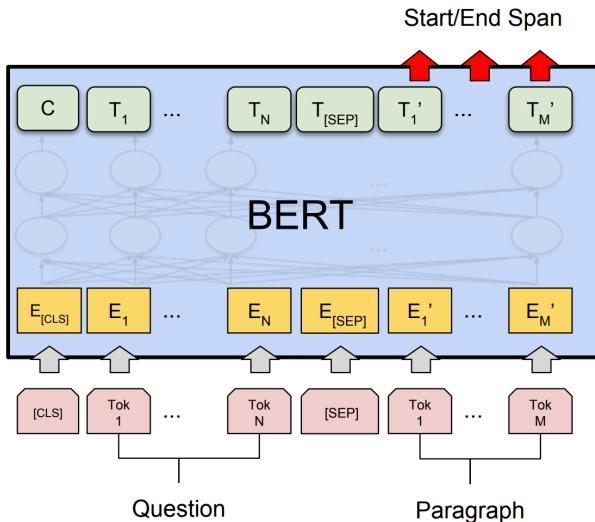


**Single sentence classification**  
Sentiment analysis, spam detection, etc.



**Pair of sentences classification**  
Entailment, paraphrase detection, etc.

# BERT – Fine-tuning for Machine Reading



(c) Question Answering Tasks:  
SQuAD v1.1

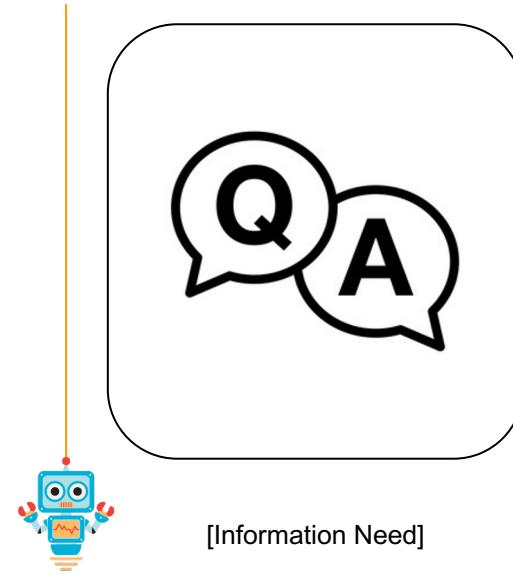
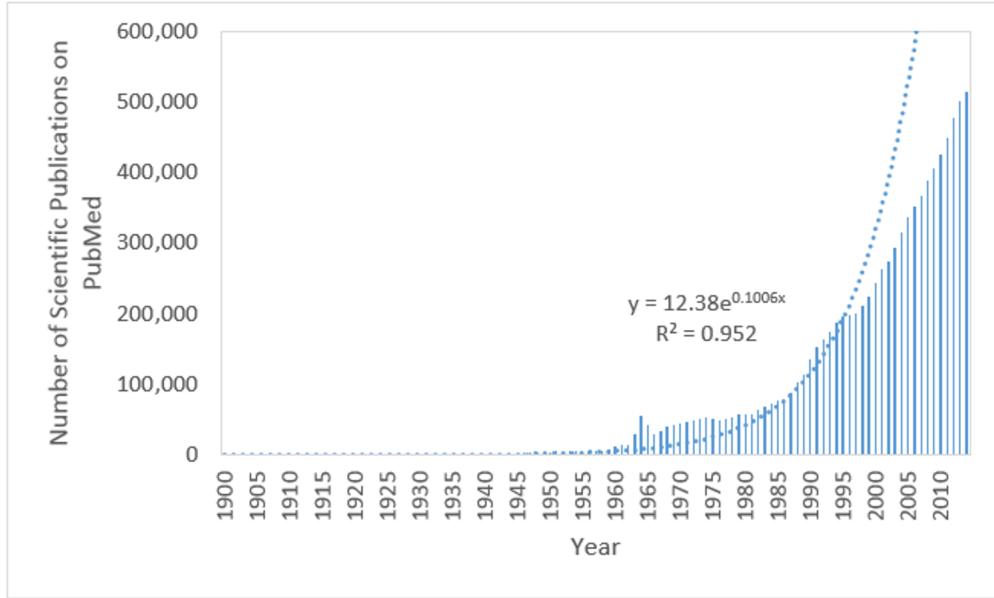
System	Dev		Test	
	EM	F1	EM	F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BiDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT <sub>BASE</sub> (Single)	80.8	88.5	-	-
BERT <sub>LARGE</sub> (Single)	84.1	90.9	-	-
BERT <sub>LARGE</sub> (Ensemble)	85.8	91.8	-	-
BERT <sub>LARGE</sub> (Sgl.+TriviaQA)	<b>84.2</b>	<b>91.1</b>	<b>85.1</b>	<b>91.8</b>
BERT <sub>LARGE</sub> (Ens.+TriviaQA)	<b>86.2</b>	<b>92.2</b>	<b>87.4</b>	<b>93.2</b>

# Open Domain Question Answering

---

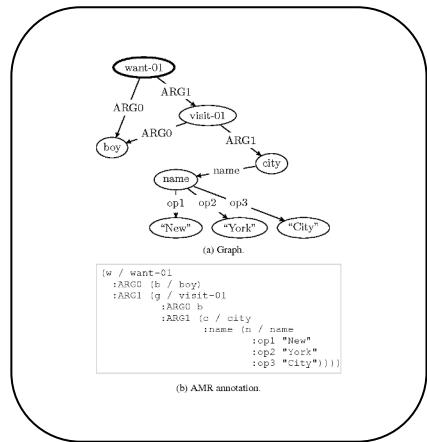
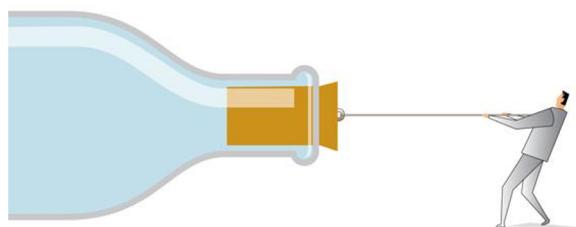
Is Machine Reading actually useful?

# Motivation 1: Information Overload



# Motivation 2: The Knowledge Acquisition Bottleneck

“The problem of knowledge acquisition is the critical bottleneck problem in artificial intelligence.”  
E. A. Feigenbaum 1984



[Meaning]



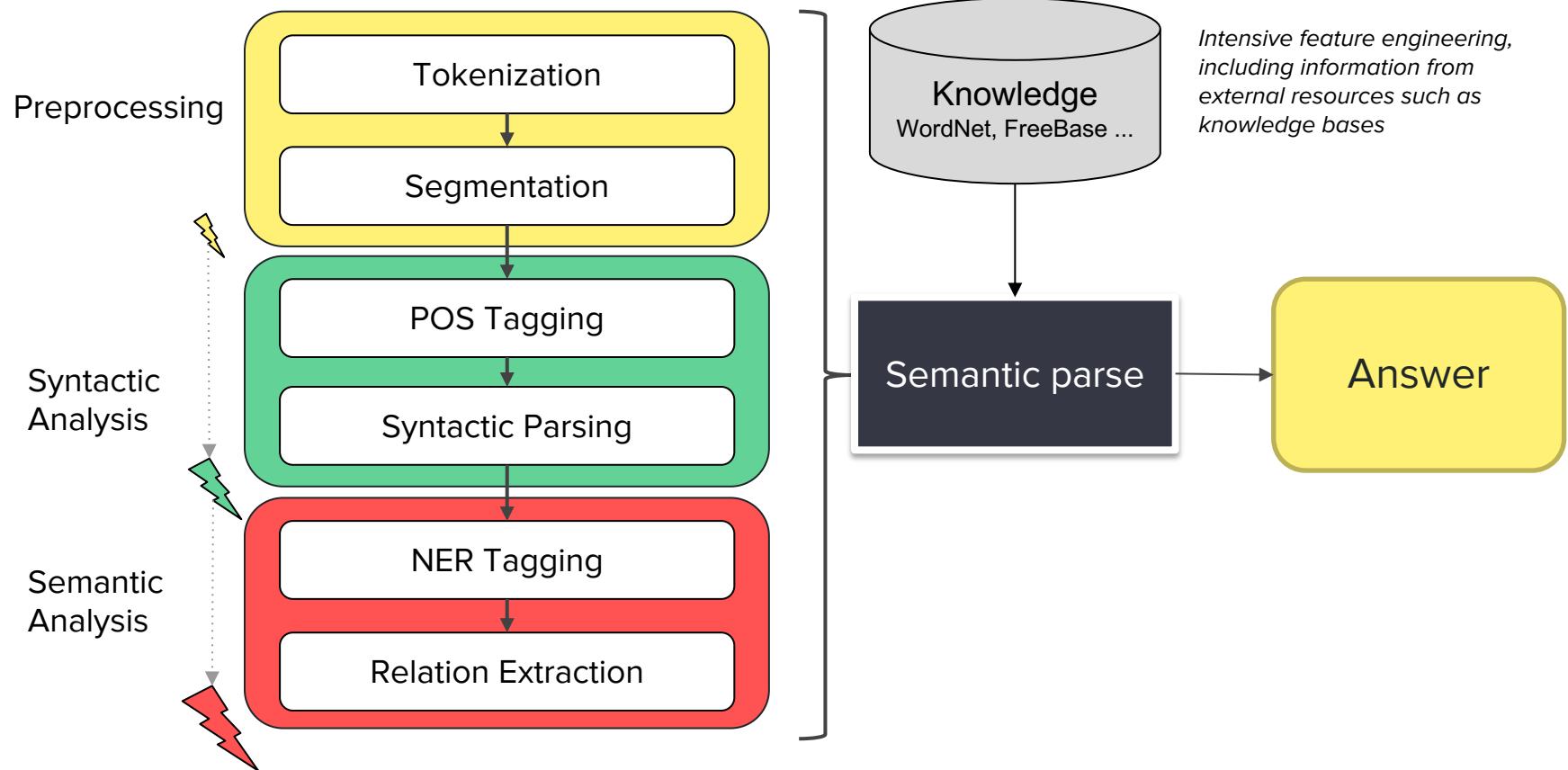
uses for



# Open domain Question Answering

- Open domain QA: answer any question using very large knowledge sources
- Goes beyond Machines Reading that expects a paragraph to be given
- Open domain = question on any topic not a restricted subset
- In the following
  1. Traditional approaches using Knowledge Bases
  2. New approaches based on end-to-end Machine Reading

# “Traditional” NLP for open domain QA



# Semantic Parsing



[Text]

$$\begin{aligned} \exists x_0 \text{named}(x_0, \text{ewan}, \text{person}) \wedge \\ \exists x_1 \text{mozzarella}(x_1) \wedge \\ \exists x_2 \text{car}(x_2) \wedge \text{of}(x_2, x_0) \wedge \text{in}(x_1, x_2) \wedge \\ \exists e \text{event}(e) \wedge \text{forget}(e) \wedge \text{agent}(e, x_0) \wedge \\ \text{patient}(e, x_1) \end{aligned}$$

[Meaning]



[Information Need]

Semantic parses are logical forms in PROLOG, SQL, SPARQL, etc.

# Knowledge Bases

- KB: structured repository of knowledge (usually relational DB)
- Goal: encode knowledge so that it can be queried by semantic parses efficiently
- Scale can be huge: billions of facts, millions of entities
- KB can be generic or specific
- Examples: Cyc, WikiData, DBPedia, Google KG, GeneOntology, IMDB, etc.

- Key challenge is their construction!

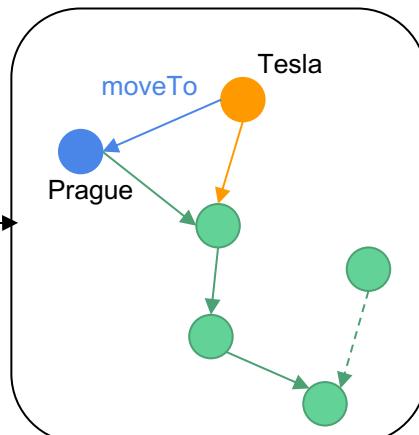
- Manually: Crowdsourcing, paid experts
- Automatically: Information extraction or Automatic KB Construction



# Automatic Knowledge Base Construction

In January 1880, two of Tesla's uncles put together enough money to help him leave Gospic for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.

[Text]

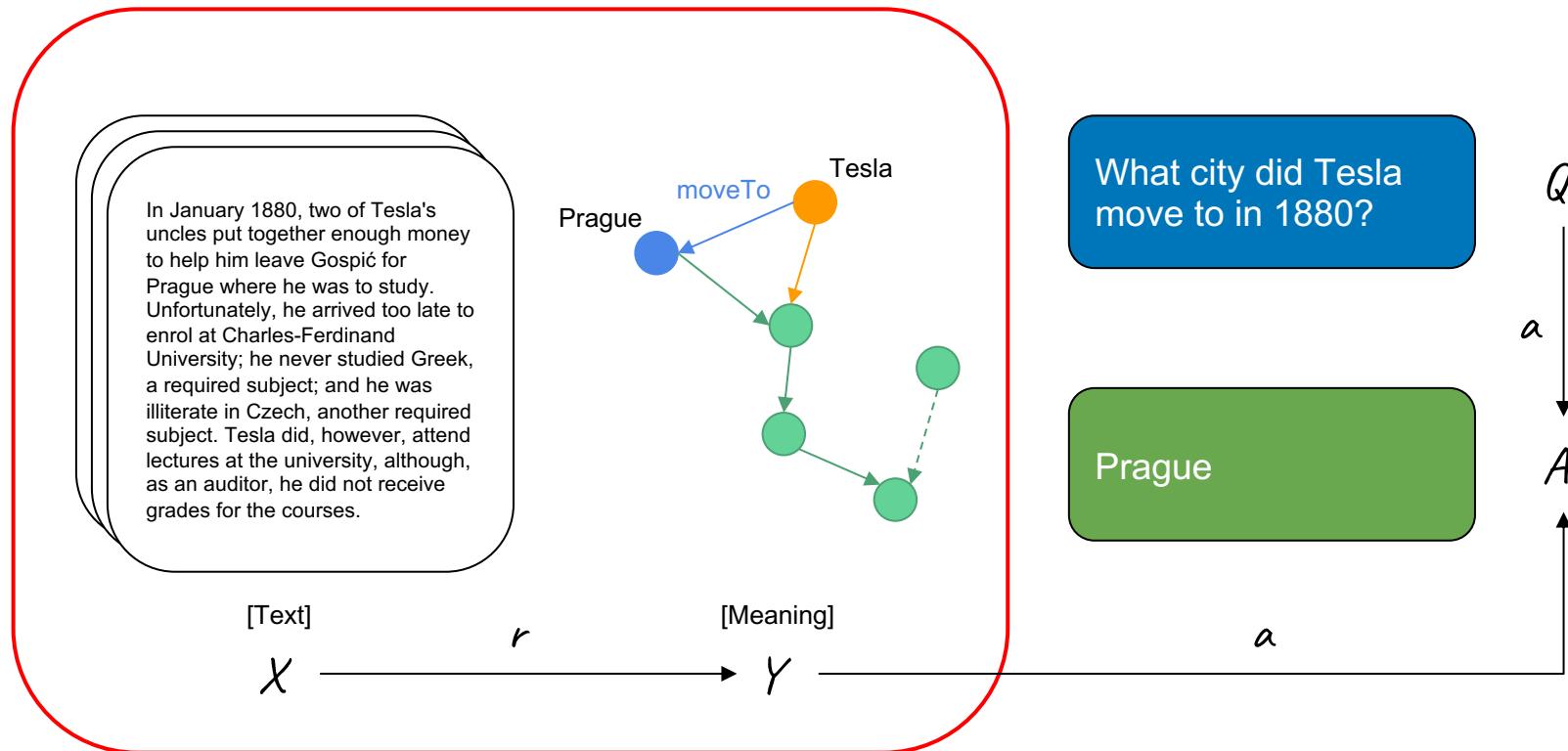


[Meaning]



[Information Need]

# Knowledge Graph Construction

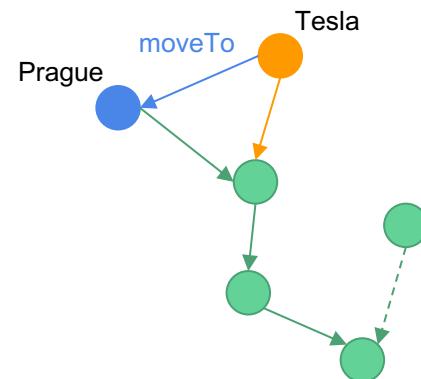


# Knowledge Graph Construction

In January 1880, two of Tesla's uncles put together enough money to help him leave Gospic for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.

[Text]

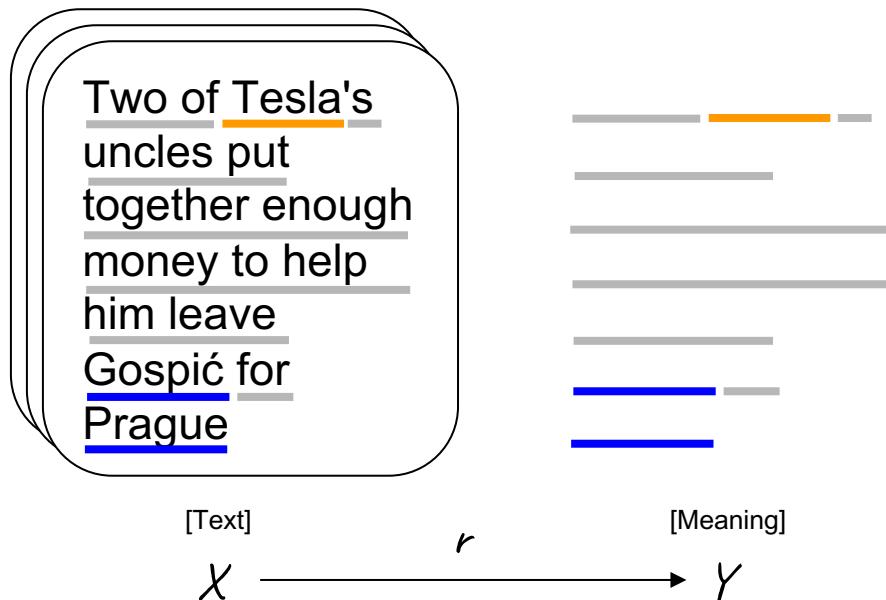
X



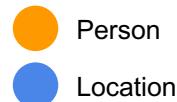
[Meaning]

Y

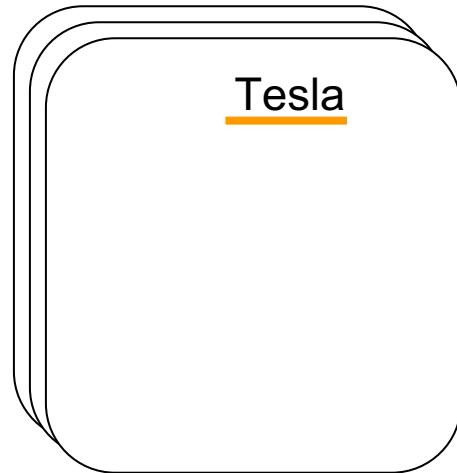
# Entity Extraction



- Linear Chain CRF
- Bi-directional RNNs
- Hybrid RNN & CRFs



# Challenge: Ambiguity

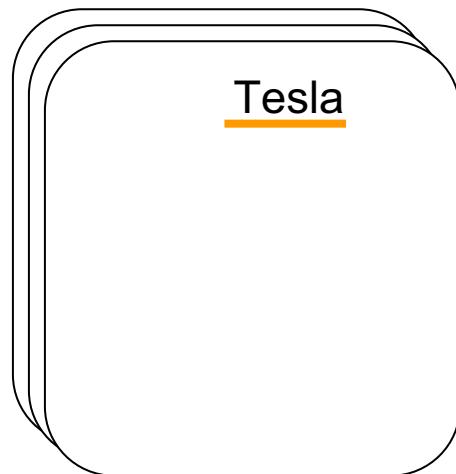


● Person?

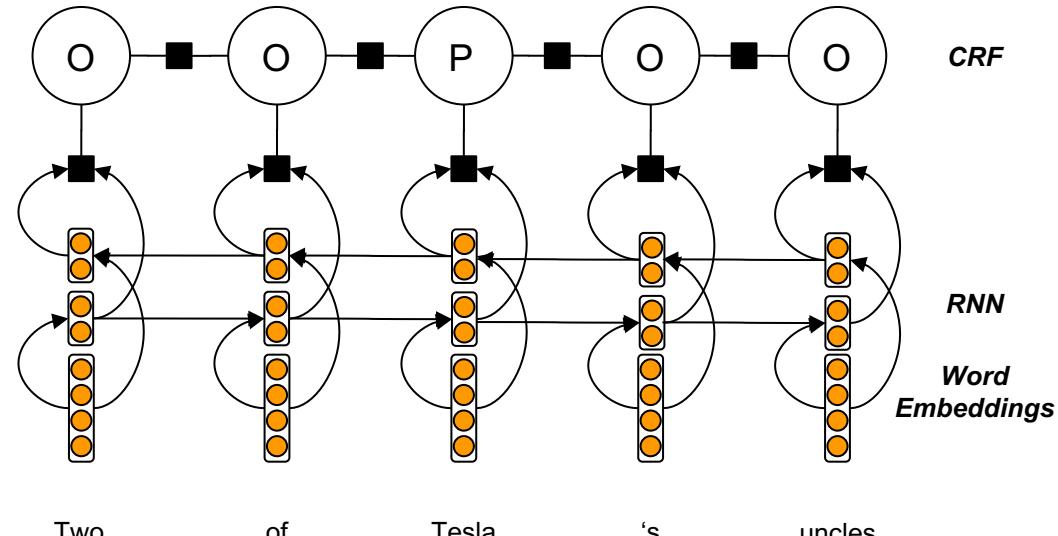
● Brand?

# Conditional Random Fields with RNN Potentials

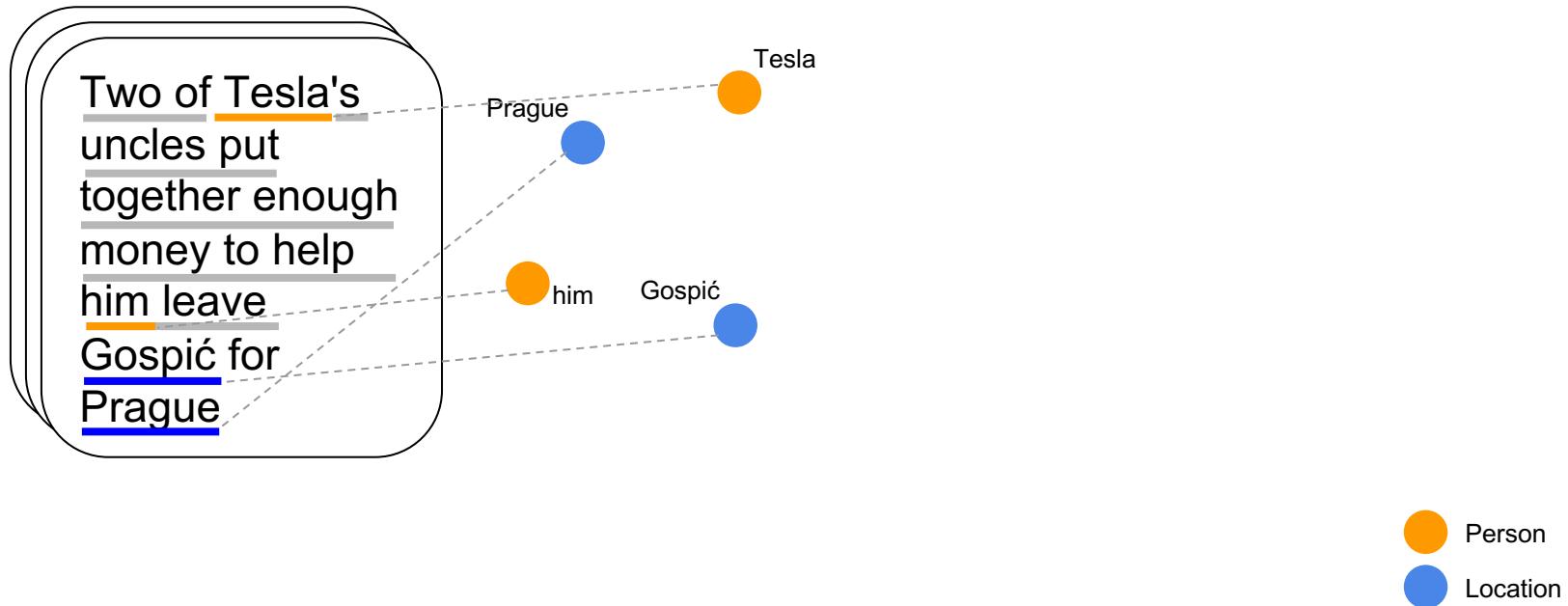
Huang et al., 2015



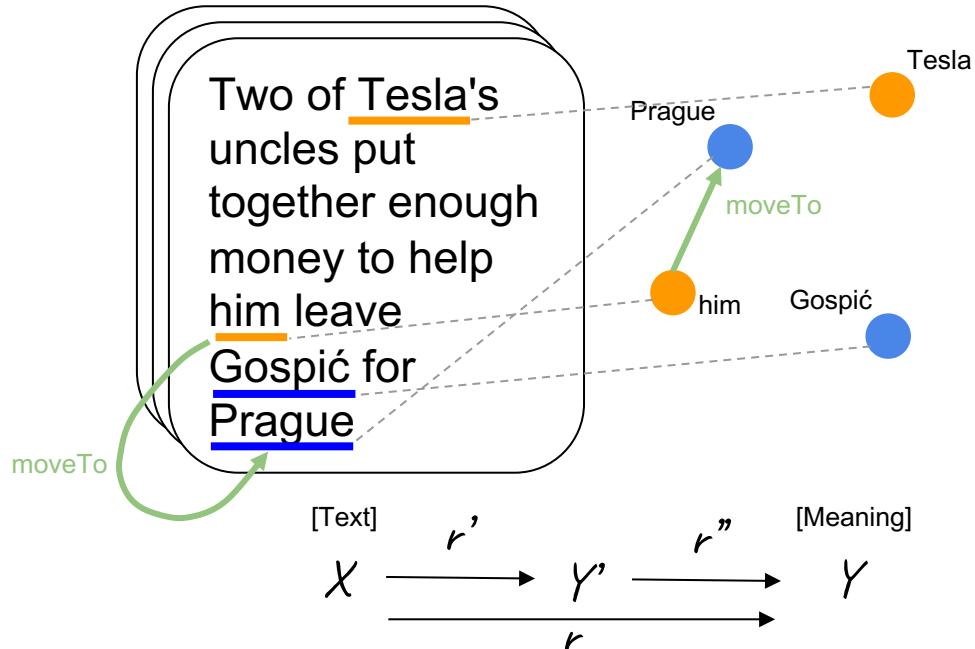
- Person?
- Brand?



# Instantiate Nodes



# Relation Extraction



- Neural Classification
- Distant Supervision

# Challenge: Variation

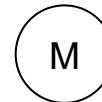
Two of Tesla's uncles put together enough money to help **him leave Gospic for Prague**

Two of Tesla's uncles put together enough money to help **him move to Prague**

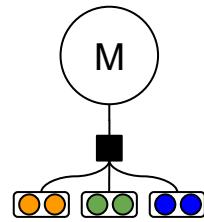
Two of Tesla's uncles put together enough money to help **him settle in Prague**

# Relation Classification

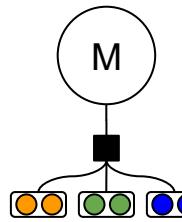
[Current SOTA neural RE model]



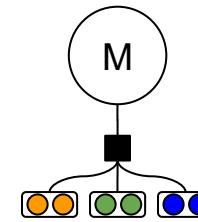
(Tesla, moveTo, Prague)



him leave  
Gospic for  
Prague



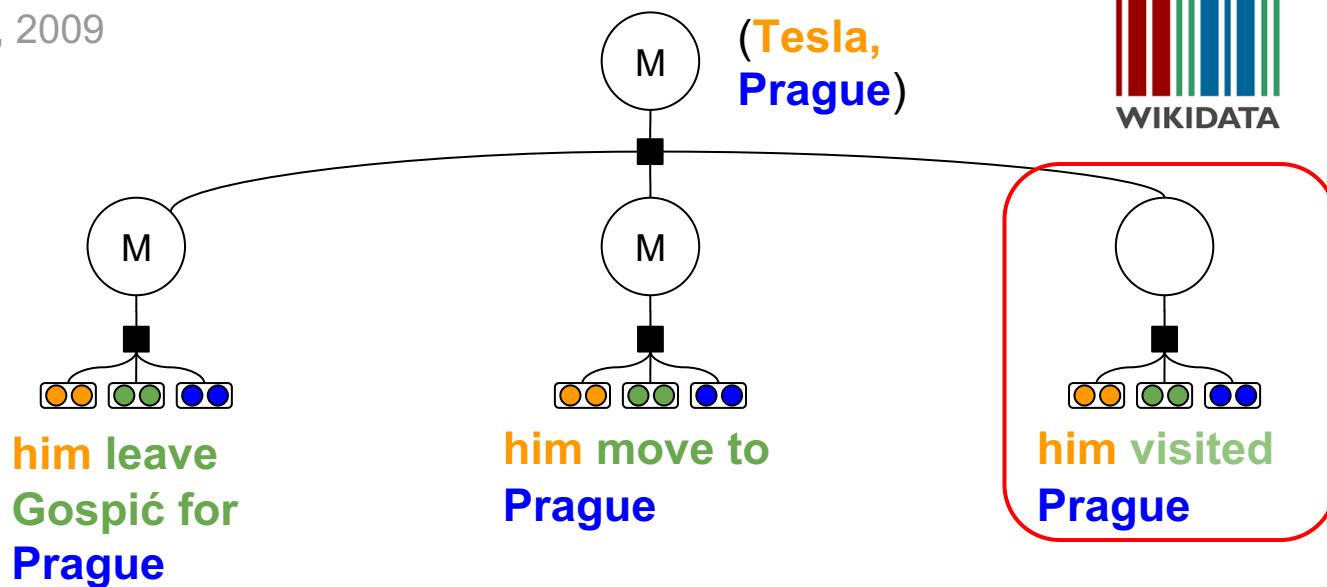
him move to  
Prague



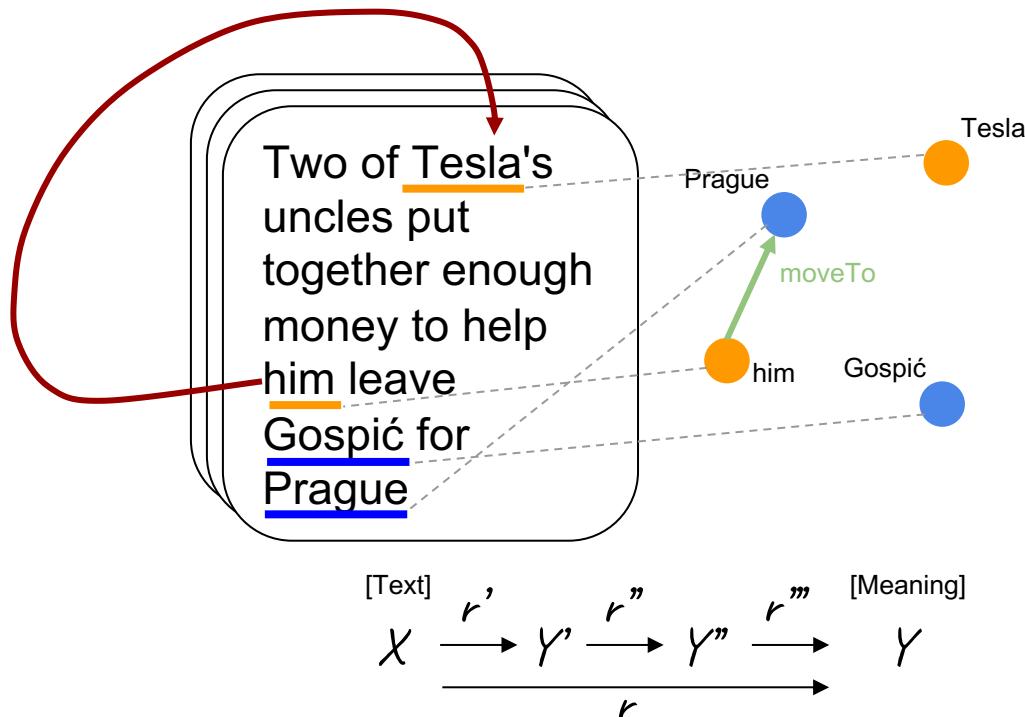
him settle in  
Prague

# Distant Supervision & Multiple Instance Learning

Mintz et al., 2009

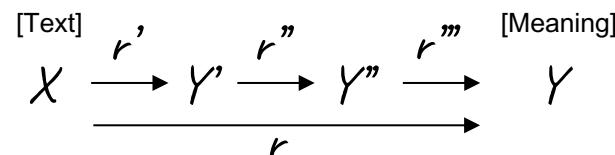
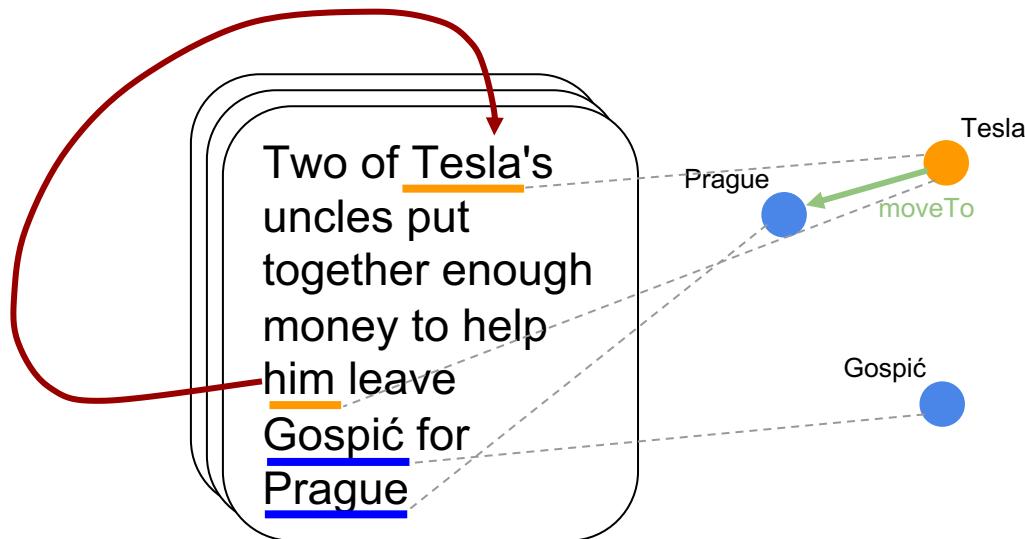


# Coreference Resolution



- Neural Classification
- Latent Variables

# Collapsing Nodes



# Challenge: Common Sense

Two of Tesla's uncles put together enough money to help him leave Gospic for Prague

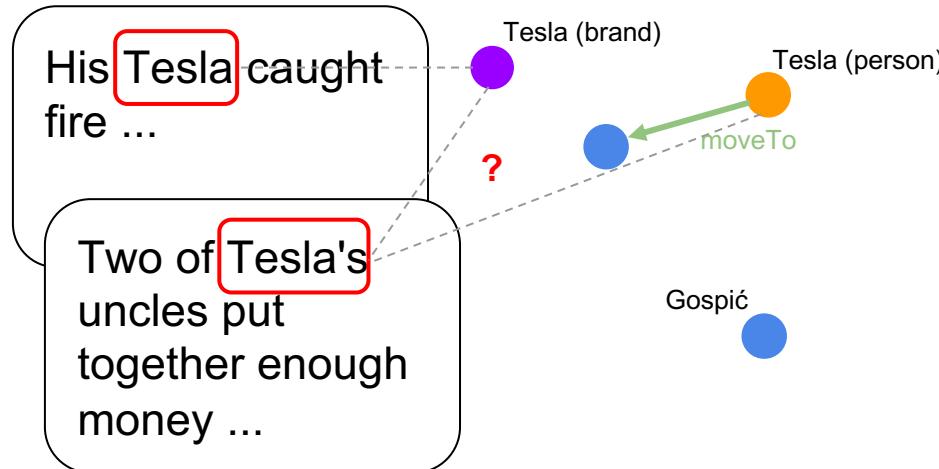
Surface

The trophy would not fit in the brown suitcase because it was too *big*.

Common Sense

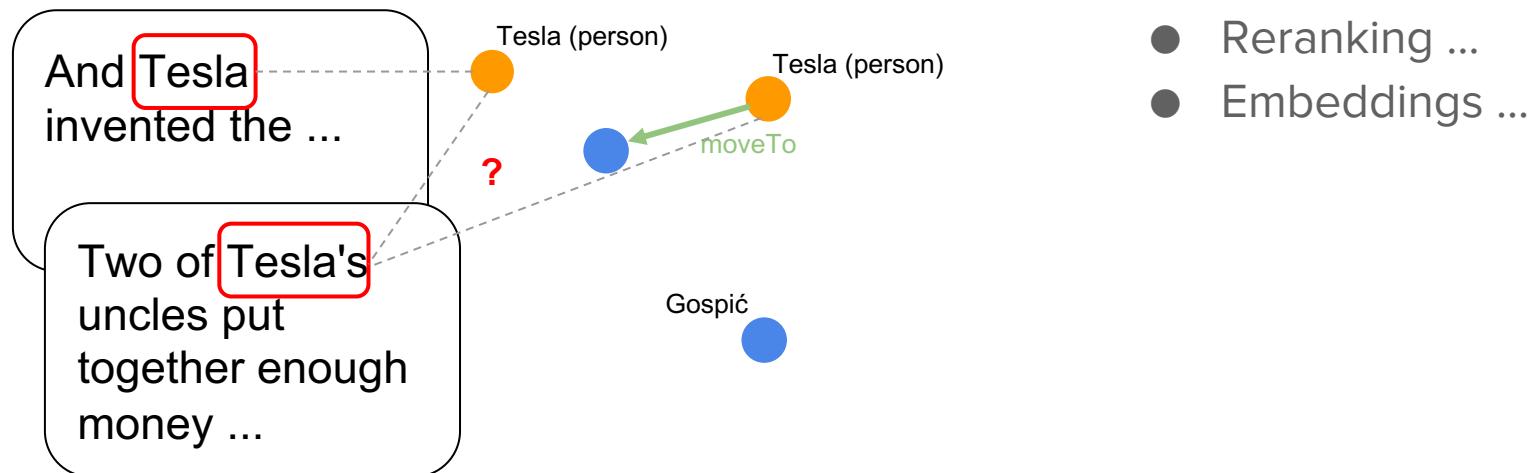
The trophy would not fit in the brown suitcase because it was too *small*.

# Entity Linking

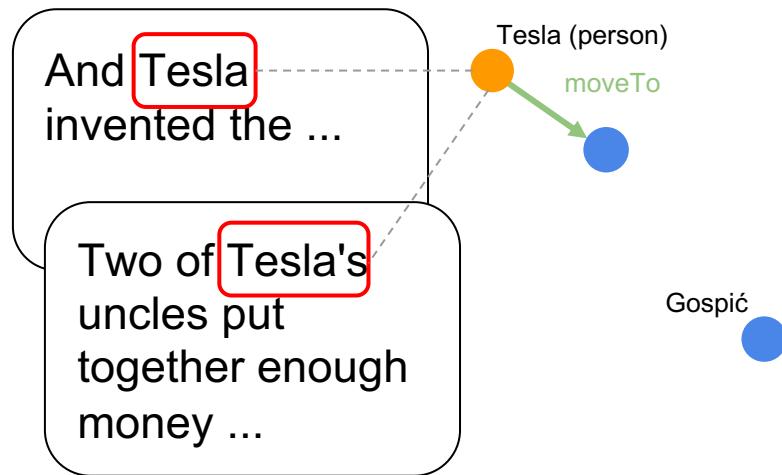


- Reranking ...
- Embeddings ...

# Entity Linking



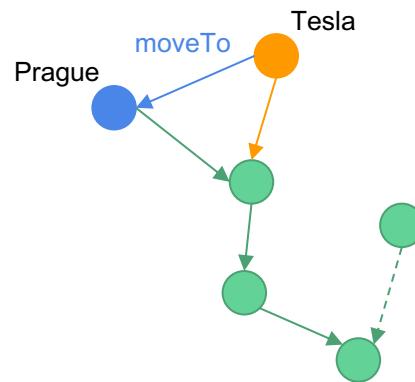
# Collapsing



$$\frac{[Text] \quad X \xrightarrow{r'} Y' \xrightarrow{r''} Y'' \xrightarrow{r'''} Y''' \xrightarrow{r''''} Y \quad [Meaning]}{r}$$

# Strengths

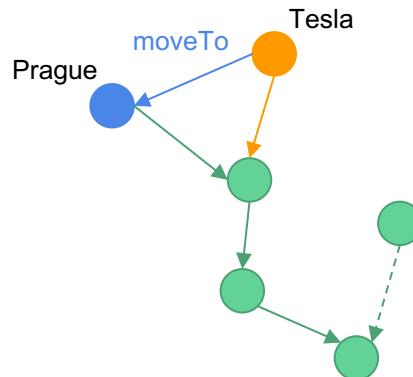
In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.



- Supports Reasoning
- Fast access
- Generalisation
- Interpretable
- Existing KBs can serve as supervision signal!

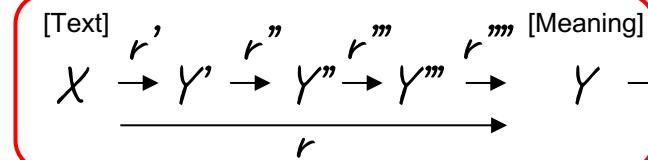
# Weakness: Cascading errors

In January 1880, two of Tesla's uncles put together enough money to help him leave Gospic for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.



What city did Tesla move to in 1880?

Prague

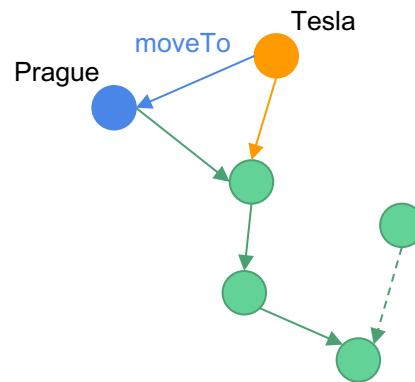


a



# Weakness: Cascading errors

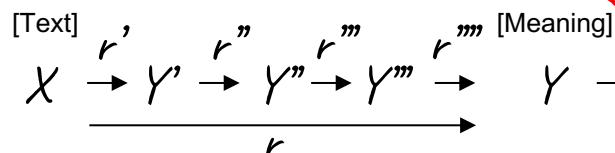
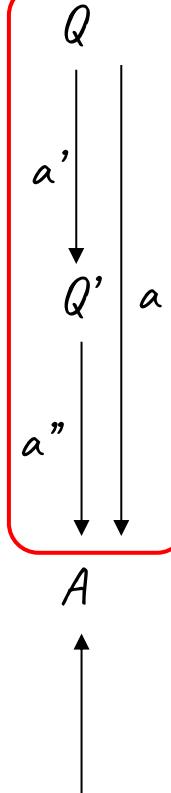
In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.



What city did Tesla move to in 1880?

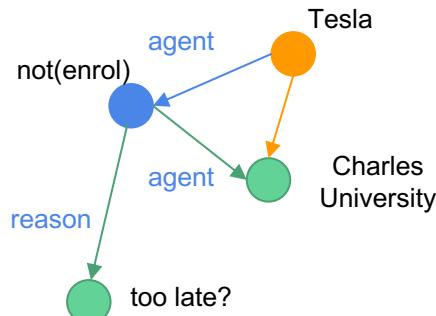
moveTo(Tesla,X)?

Prague



# Weakness: Engineering Schemas and Formalisms

Unfortunately, he arrived too late to enrol at Charles University



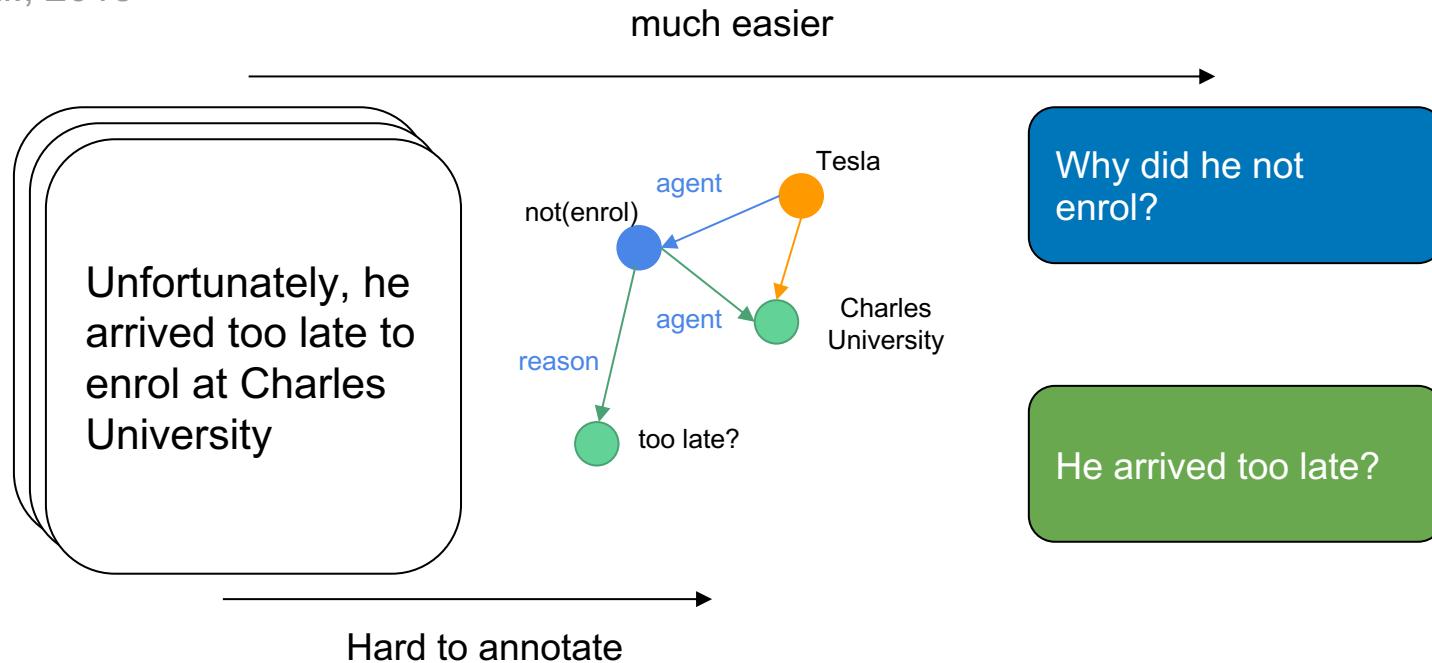
Why did he not enrol?

He arrived too late?

getting this right is hard

# Weakness: Annotation

He et al., 2015



# Structured Representations

- Advantages

- Fast access
- Scalable
- Interpretable
- Supports reasoning
- Universality of representations: independent of question

- Disadvantages

- Less robust to variation in language
- Cascading errors
- Schema engineering
- Annotation requires experts

# Is there another way?

In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.

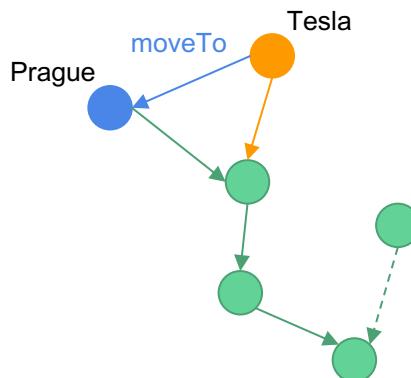
[Text]

$X$

$r$

[Meaning]

$Y$



What city did Tesla move to in 1880?

Prague

$Q$

$a$

$A$

$a$

# Omitting Symbolic Meaning Representations !!

In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for Prague where he was to study. Unfortunately, he arrived too late to enrol at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses.

[Text]

X

—

What city did Tesla move to in 1880?

Prague

a

Q

a

A

# Machine Reading AT SCALE

A **machine** processes a (very) large collection of texts to satisfy an **information need**

# Machine Reading



[Text]

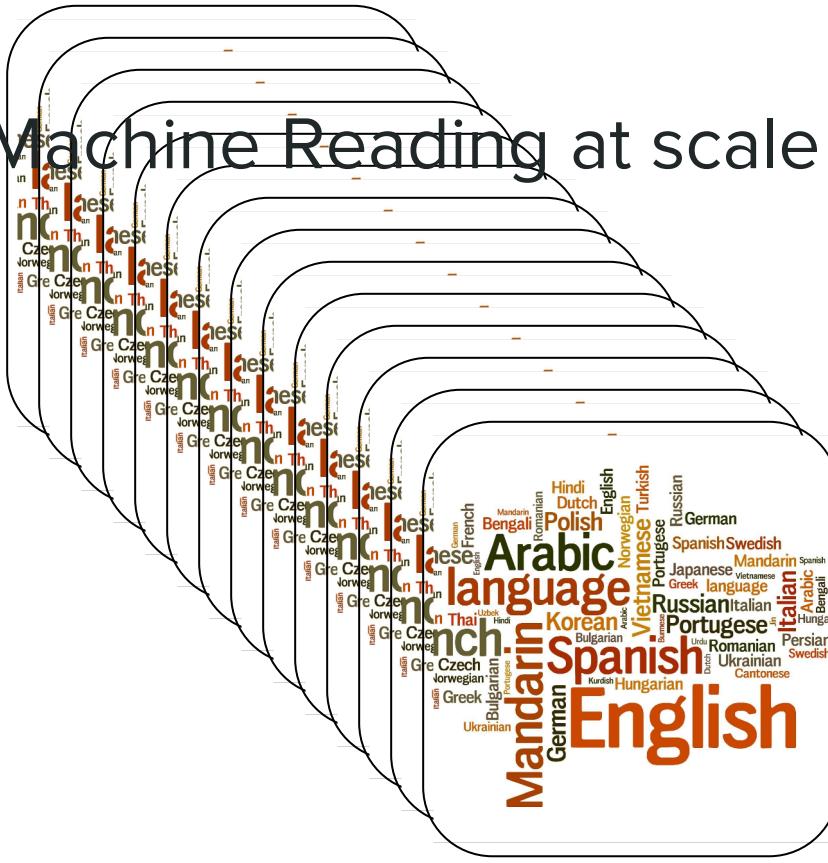


## **uses for**



## [Information Need]

# Machine Reading at scale



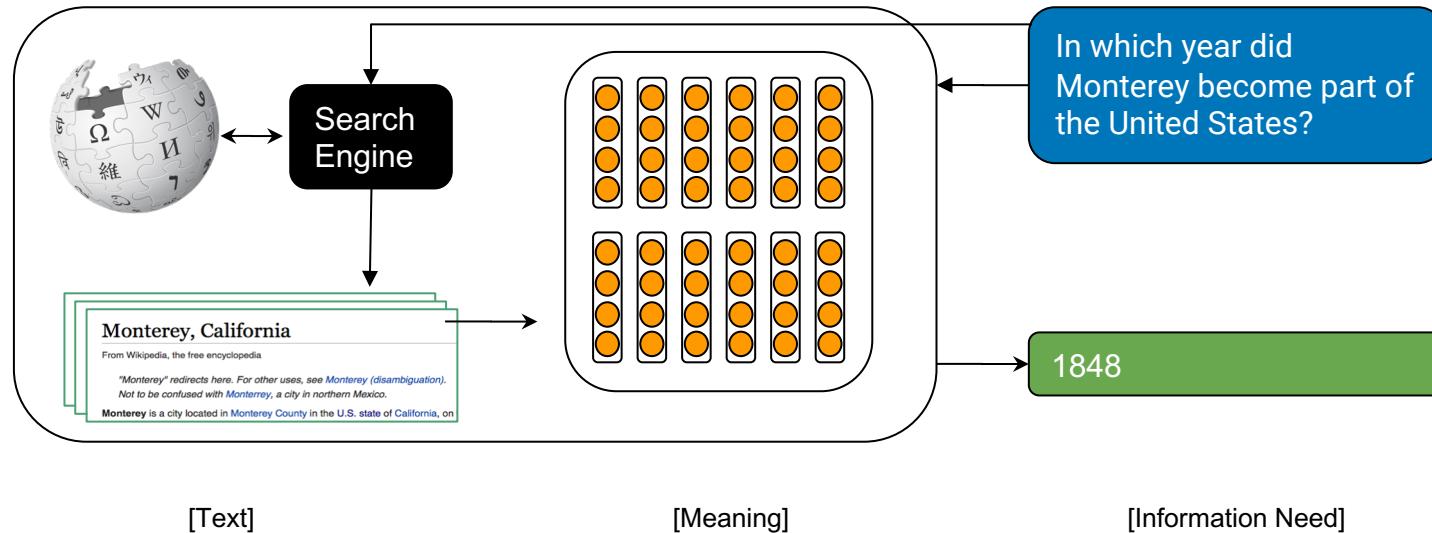
uses for



[Information Need]

# Typical Machine Reading at Scale System

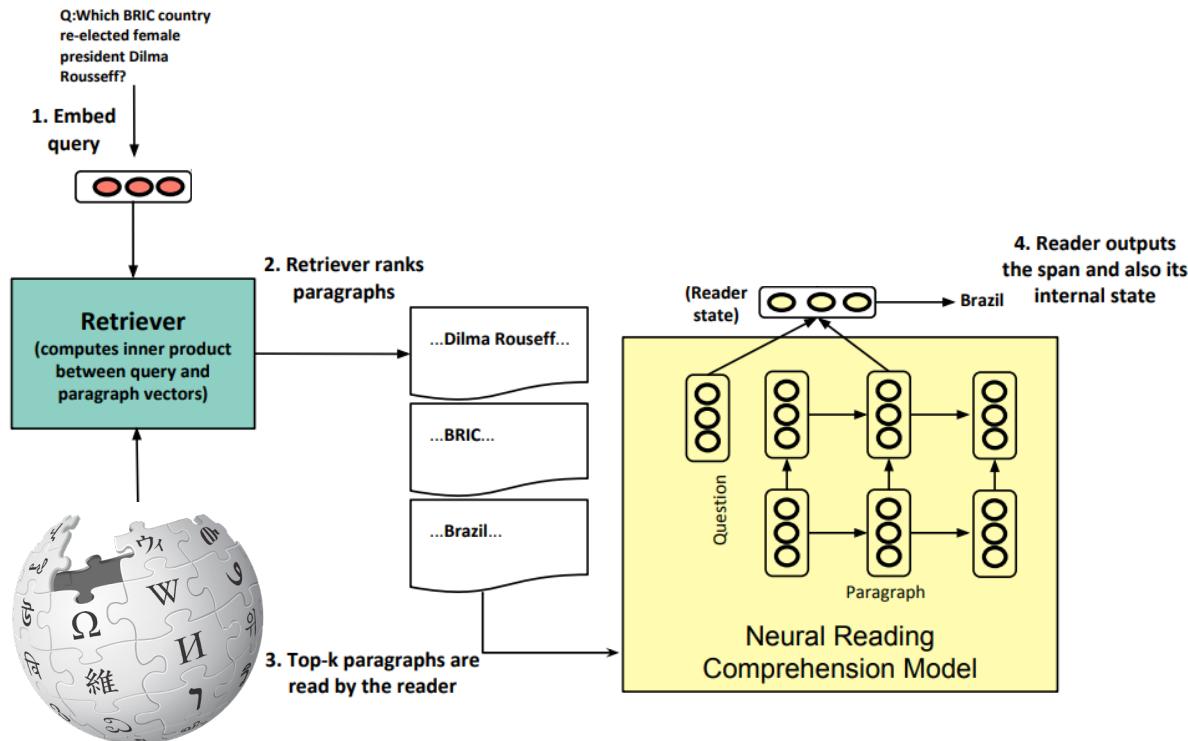
Dr.QA Chen et al., 2017



No way to recover if the search engine is wrong!

# Current best: Multi-Step Retriever-Reader

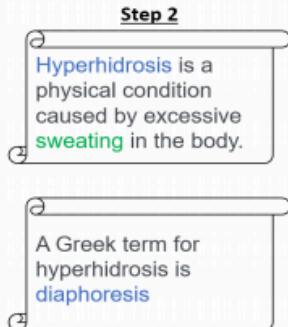
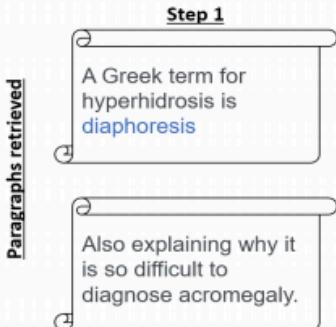
Das et al., 2019



# Current best: Multi-Step Retriever-Reader

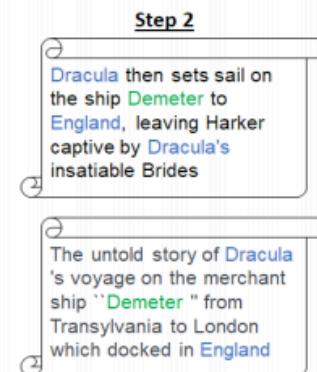
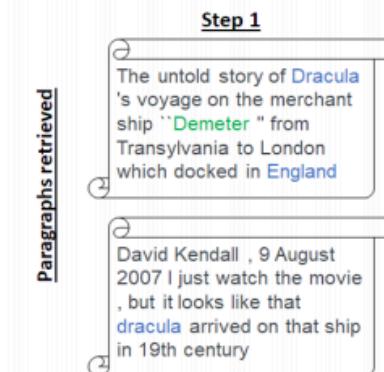
Das et al., 2019

Query: "Diaphoresis" is a medical term for what condition?



Answer: sweating

Query: What is name of the ship on which Dracula arrived in England in 1897?

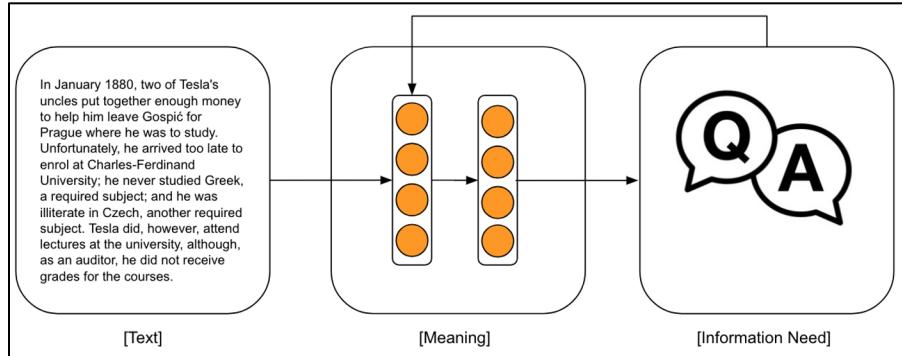
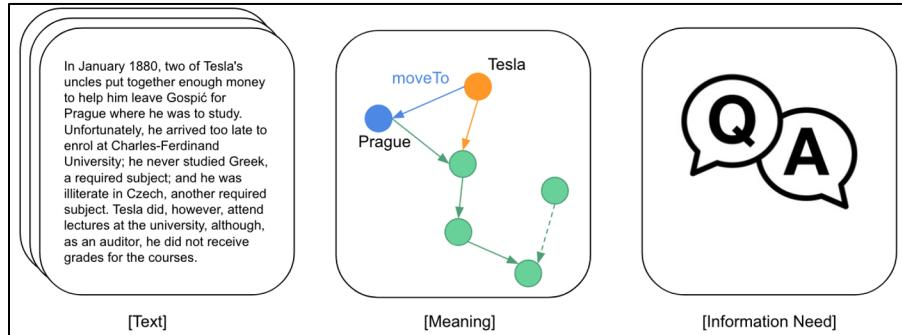


Answer: demeter

Between 40 and 60% of correct responses (for rather simple questions)

# A Paradigm Shift

- Symbolic Meaning Representations  
→ Latent Vector Representations
- Feature Engineering & Domain Expertise  
→ Architecture Engineering & ML/DL Expertise



# Pros and cons

End-to-end models	Symbolic systems
<p><i>Neural Networks</i></p> <ul style="list-style-type: none"><li>• Scale to very large datasets</li><li>• Can be used by non domain experts</li><li>• Robust to noise and ambiguity in data</li><li>• Game changers in multiple applications</li><li>• Very data hungry (mostly supervised data)</li><li>• Can't learn easily new tasks from old ones</li><li>• Not interpretable</li><li>• Relatively simple reasoning</li></ul>	<p><i>KBs, Inductive Logic Programming, etc.</i></p> <ul style="list-style-type: none"><li>• Small scale conditions</li><li>• Require heavy expert knowledge</li><li>• Very brittle with noisy, ambiguous data</li><li>• Limited applicative success</li></ul> <p>Great research opportunities!</p>

# Current Challenge: Reconciling Conflicting Information

*So how much does the UK pay to the EU per week?*

“Once we have settled our accounts, we will take back control of roughly **£350m** per week.” *Boris Johnson*

“We are not giving £20bn a year or £350m a week to Brussels - Britain pays **£276m** a week to the EU budget because of the rebate.” *BBC Reality Check*

“...When those are taken into account the figure is **£250m.**” *Independent*

?

Trust into source, timeline, ...

# Conclusion

- We've seen 2 approaches for building system to answer any question
- Most deployed systems still rely on traditional pipelines for the most part (+ some DL here and there)
- Why? **Scale, reliability, interpretability**
- Open questions:
  - All shortcomings of Machine Reading → Open domain QA. Need to solve them
  - Will pretrained contextual embeddings change everything forever?
  - Can we combine both symbolic and end-to-end approaches?